

Code Generation: Results Achieved with Zero-Shot prompting when not asking for a Rationale

Table 1: Number (#) and percentage (%) of instances for which the LLMs did not manage to output a valid judgment.

| LLM | Code generation | | | |
|---------------------|-----------------|--------|--------|-------|
| | Java | | Python | |
| | # | % | # | % |
| DeepSeek Coder 1.3B | 379 | 26.98% | 31 | 2.42% |
| DeepSeek Coder 6.7B | 0 | 0.00% | 3 | 0.23% |
| DeepSeek Coder 33B | 3 | 0.21% | 3 | 0.23% |
| CodeLlama 7B | 1 | 0.07% | 61 | 4.76% |
| CodeLlama 13B | 3 | 0.21% | 44 | 3.43% |
| CodeLlama 34B | 3 | 0.21% | 9 | 0.70% |
| GPT-3.5-turbo | 0 | 0.00% | 0 | 0.00% |
| GPT-4-turbo | 0 | 0.00% | 0 | 0.00% |

| | DSC 1.3B | DSC 6.7B | DSC 33B | CL 7B | CL 13B | CL 34B | GPT 3.5 | GPT 4 | Human Written | Own vs LLMs | Own vs LLMs \ F | Own vs Human |
|-----------------|-------------|-------------|------------|----------|-----------|-----------|------------|----------|------------------|----------------|--------------------|-----------------|
| DSC 1.3B | 0.64 | 0.58 | 0.47 | 0.57 | 0.47 | 0.54 | 0.49 | 0.57 | -0.11 | (N) | (N) | *** (L) |
| DSC 6.7B | 0.78 | 0.76 | 0.65 | 0.76 | 0.63 | 0.68 | 0.65 | 0.64 | -0.01 | *(N) | *(N) | *** (L) |
| DSC 33B | 0.35 | 0.42 | 0.30 | 0.39 | 0.32 | 0.35 | 0.37 | 0.40 | -0.60 | (N) | (N) | *** (L) |
| CL 7B | -0.12 | -0.18 | -0.24 | -0.13 | -0.28 | -0.20 | -0.14 | -0.24 | -0.93 | (N) | (N) | *** (L) |
| CL 13B | 0.62 | 0.56 | 0.43 | 0.61 | 0.46 | 0.49 | 0.61 | 0.57 | -0.09 | *(N) | *(N) | *** (L) |
| CL 34B | 0.23 | 0.22 | 0.12 | 0.21 | 0.17 | 0.13 | 0.32 | 0.30 | -0.52 | (N) | *(N) | *** (L) |
| GPT-3.5 | 0.53 | 0.62 | 0.48 | 0.59 | 0.51 | 0.55 | 0.57 | 0.55 | -0.23 | (N) | (N) | *** (L) |
| GPT-4 | 0.34 | 0.44 | 0.34 | 0.38 | 0.35 | 0.33 | 0.49 | 0.52 | -0.42 | ** (N) | ** (N) | *** (L) |
| Average (all) | 0.42 | 0.45 | 0.34 | 0.44 | 0.36 | 0.37 | 0.47 | 0.47 | -0.37 | - | - | - |
| Average (large) | 0.42 | 0.43 | 0.32 | 0.42 | 0.33 | 0.36 | 0.42 | 0.41 | -0.36 | - | - | - |

Adjusted p -values: * <0.05, ** <0.01, *** <0.001. Cliff delta: N=Negligible, S=Small, M=Medium, L=Large

Table 2: **(Java)** Average of differences between the LLM judgments (0 or 1) and the ground truth (*i.e.*, 1 if the method passes the test and 0 otherwise). Last three columns report adj. p -value and effect size when comparing the judgements each LLM gave to functions it generated against those it gave when judging functions (i) generated by all other LLMs, (ii) generated by all other LLMs but those belonging to the same family, and (iii) written by humans.

| | DSC 1.3B | DSC 6.7B | DSC 33B | CL 7B | CL 13B | CL 34B | GPT 3.5 | GPT 4 | Human Written | Own vs LLMs | Own vs LLMs \ F | Own vs Human |
|-----------------|-------------|-------------|------------|----------|-----------|-----------|------------|----------|------------------|----------------|--------------------|-----------------|
| DSC 1.3B | 0.78 | 0.72 | 0.69 | 0.76 | 0.60 | 0.71 | 0.65 | 0.63 | -0.11 | *(N) | *(N) | *** (L) |
| DSC 6.7B | 0.81 | 0.81 | 0.79 | 0.78 | 0.68 | 0.73 | 0.72 | 0.74 | -0.03 | (N) | (N) | *** (L) |
| DSC 33B | 0.31 | 0.46 | 0.42 | 0.38 | 0.38 | 0.34 | 0.40 | 0.53 | -0.45 | (N) | (N) | *** (L) |
| CL 7B | 0.01 | -0.01 | -0.08 | 0.07 | -0.17 | -0.14 | -0.13 | -0.04 | -0.91 | ** (N) | *(N) | *** (L) |
| CL 13B | 0.59 | 0.59 | 0.62 | 0.60 | 0.52 | 0.55 | 0.43 | 0.52 | -0.23 | (N) | (N) | *** (L) |
| CL 34B | 0.47 | 0.48 | 0.52 | 0.52 | 0.46 | 0.47 | 0.45 | 0.52 | -0.30 | (N) | (N) | *** (L) |
| GPT-3.5 | 0.64 | 0.71 | 0.64 | 0.63 | 0.55 | 0.63 | 0.63 | 0.69 | -0.22 | (N) | (N) | *** (L) |
| GPT-4 | 0.38 | 0.59 | 0.58 | 0.48 | 0.43 | 0.48 | 0.58 | 0.71 | -0.36 | *** (S) | *** (S) | *** (L) |
| Average (all) | 0.48 | 0.56 | 0.56 | 0.52 | 0.47 | 0.49 | 0.50 | 0.59 | -0.31 | - | - | - |
| Average (large) | 0.50 | 0.54 | 0.52 | 0.52 | 0.43 | 0.47 | 0.47 | 0.54 | -0.33 | - | - | - |

Adjusted p -values: * <0.05, ** <0.01, *** <0.001. Cliff delta: N=Negligible, S=Small, M=Medium, L=Large

Table 3: **(Python)** Average of differences between the LLM judgments (0 or 1) and the ground truth (*i.e.*, 1 if the method passes the test and 0 otherwise). Last three columns report adj. p -value and effect size when comparing the judgements each LLM gave to functions it generated against those it gave when judging functions (i) generated by all other LLMs, (ii) generated by all other LLMs but those belonging to the same family, and (iii) written by humans.

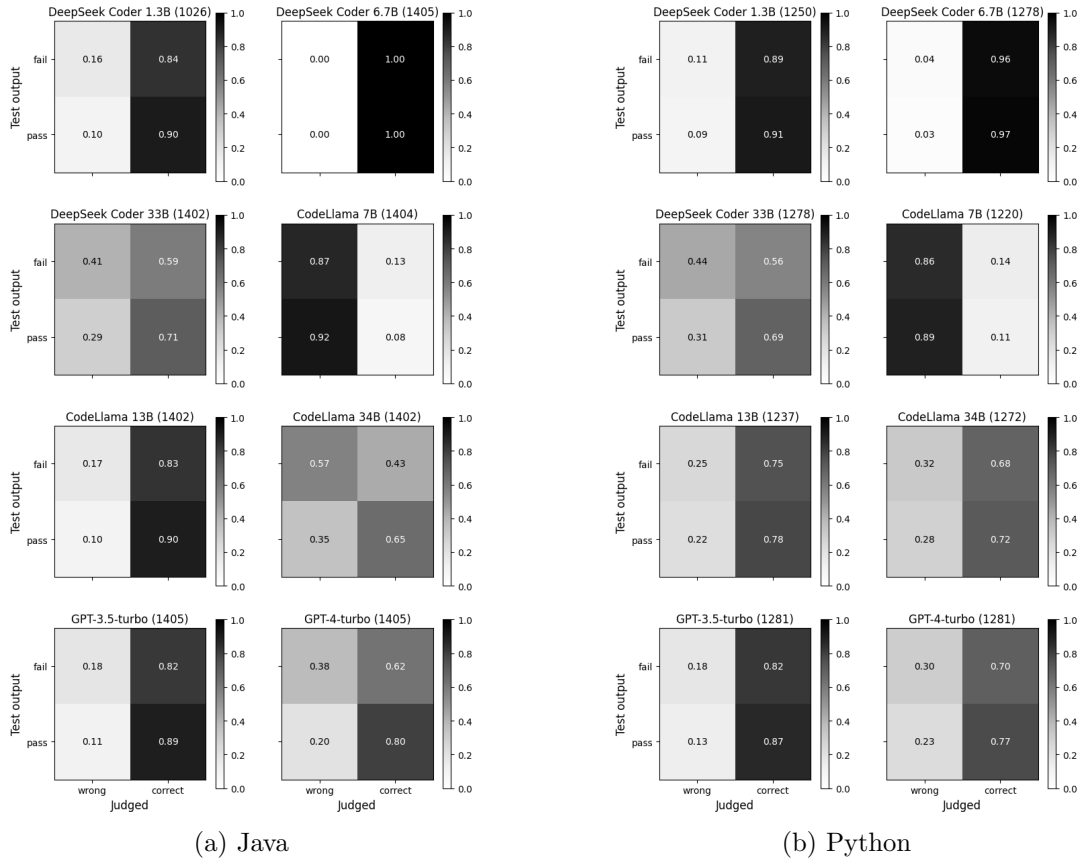


Figure 1: Code Generation: Confusion matrices for LLM's judgment.